

# Reaching Petabyte-Scale System with ZNS SSDs

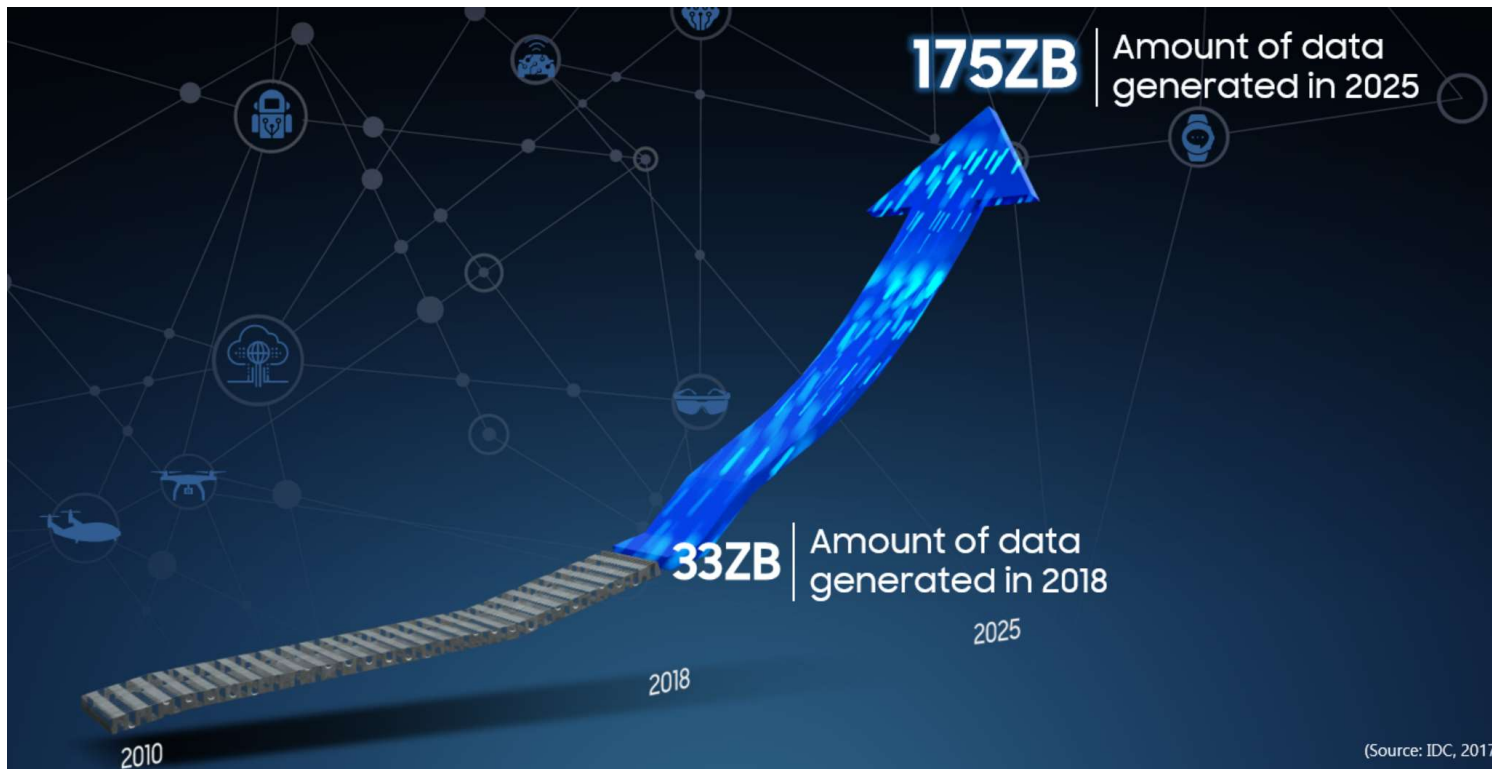
2022. 11. 01 | Wonchul Lee,  
Samsung Electronics

- The Advent of the Data Age
- Challenges
  - Zone & Data Management Overhead
  - Blast Radius
  - Initialization Time
  - Others
- Summary

# The Advent of the Data Age

Core Confidential

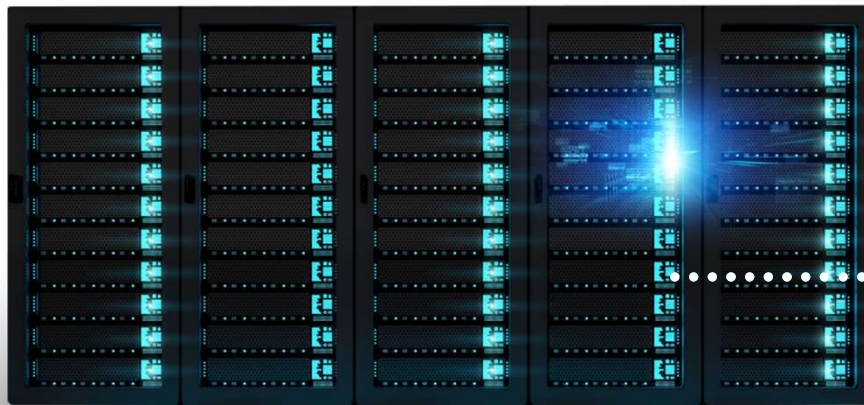
- Data is increasing at an exponential rate and could reach 175ZB by 2025.



# Petabyte SSD (PBSSD)

Core Confidential

- Disaggregated NVMe subsystem with QLC SSDs increases rack-scale space and power efficiency
- Building block for a super high-density storage system



**Exabyte scale storage system**

**1PB disaggregated storage box**

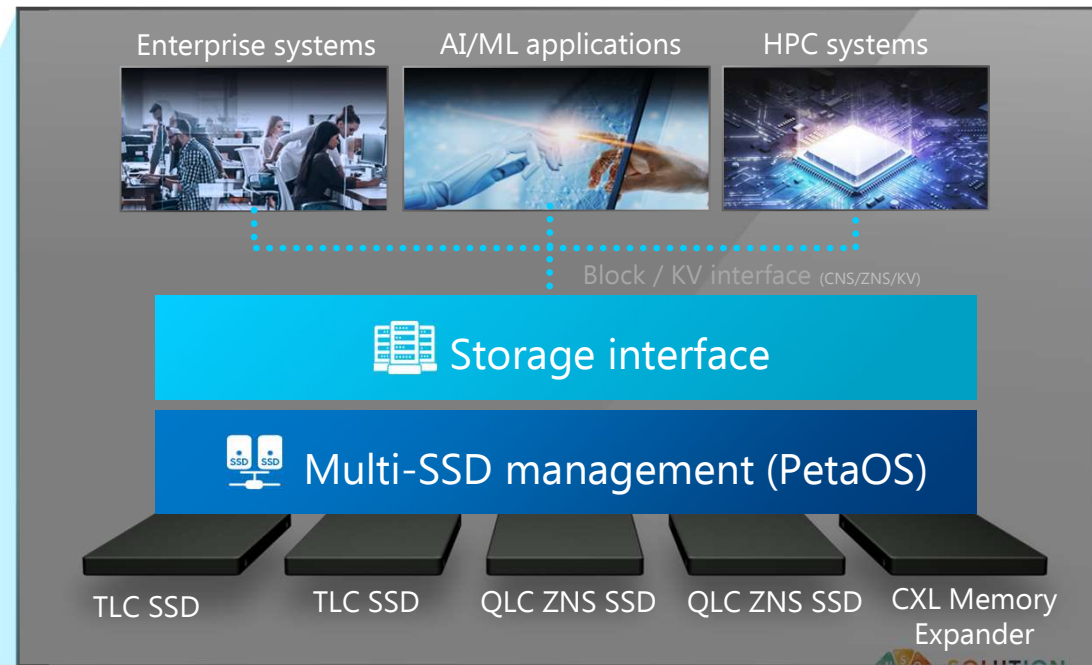
Space efficiency | Power efficiency | High performance



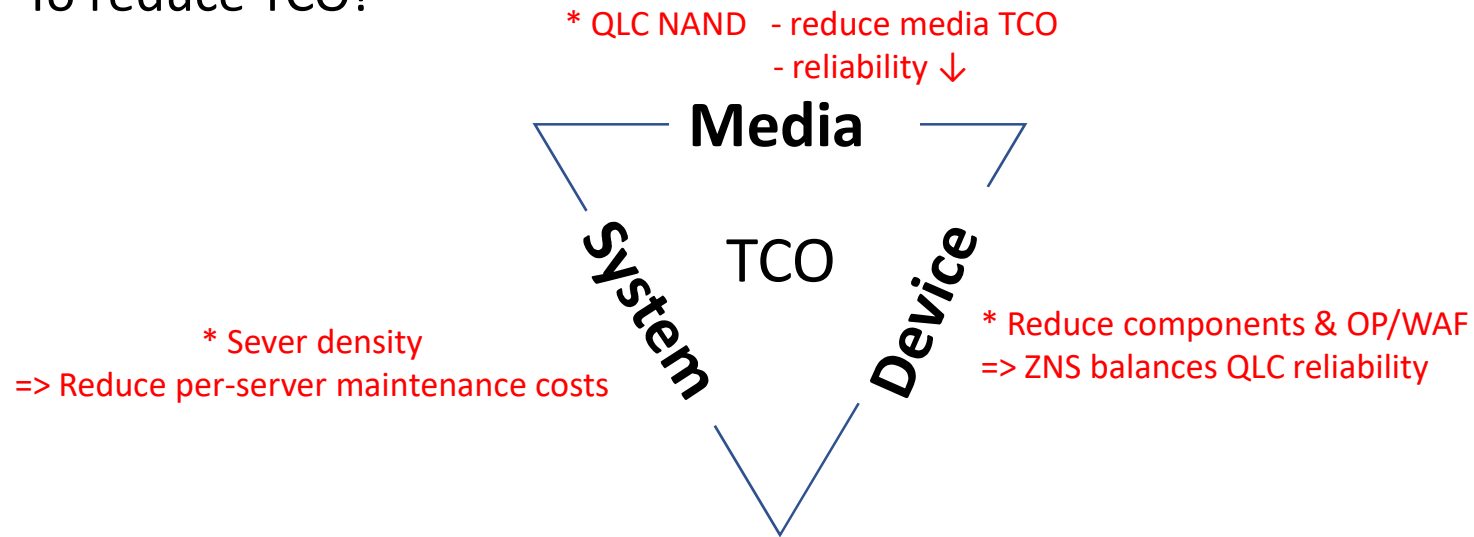
# Petabyte SSD (PBSSD)

Core Confidential

- Disaggregated NVMe subsystem with QLC SSDs increases rack-scale space and power efficiency
- Building block for a super high-density storage system



- To reduce TCO!

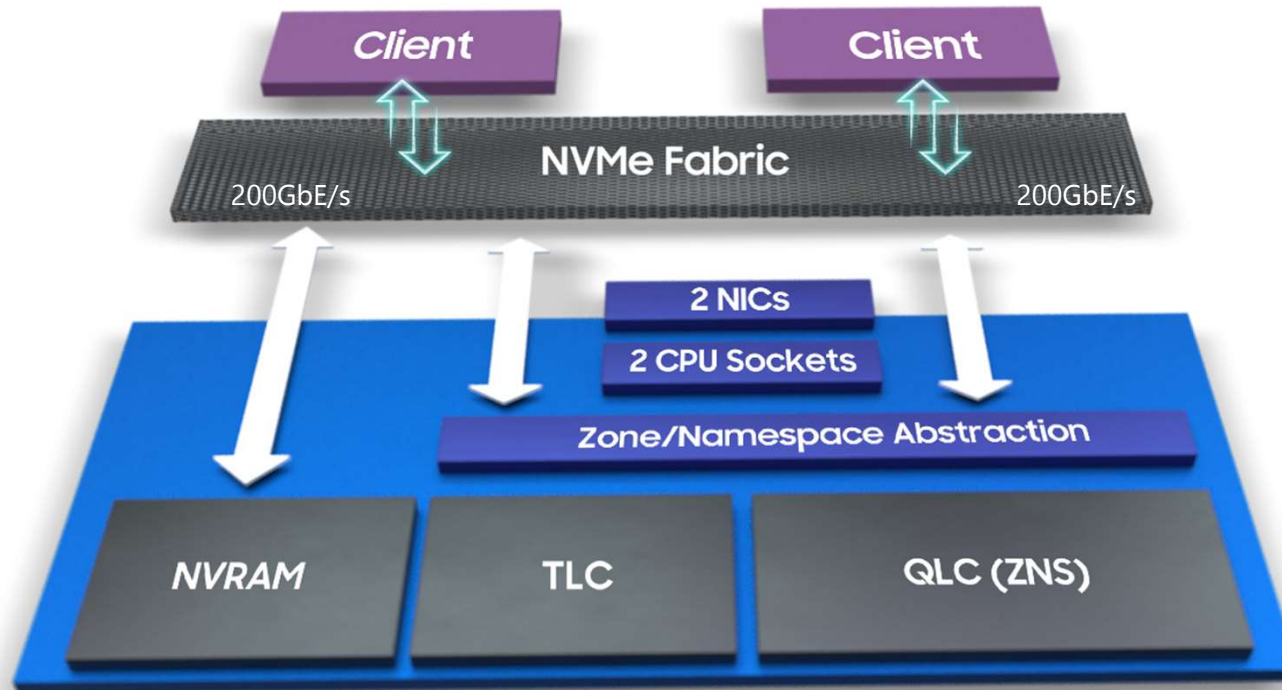


- Challenges
  - Zone & Data management O/H
  - Blast radius
  - Bring-up time

# PBSSD Reference System

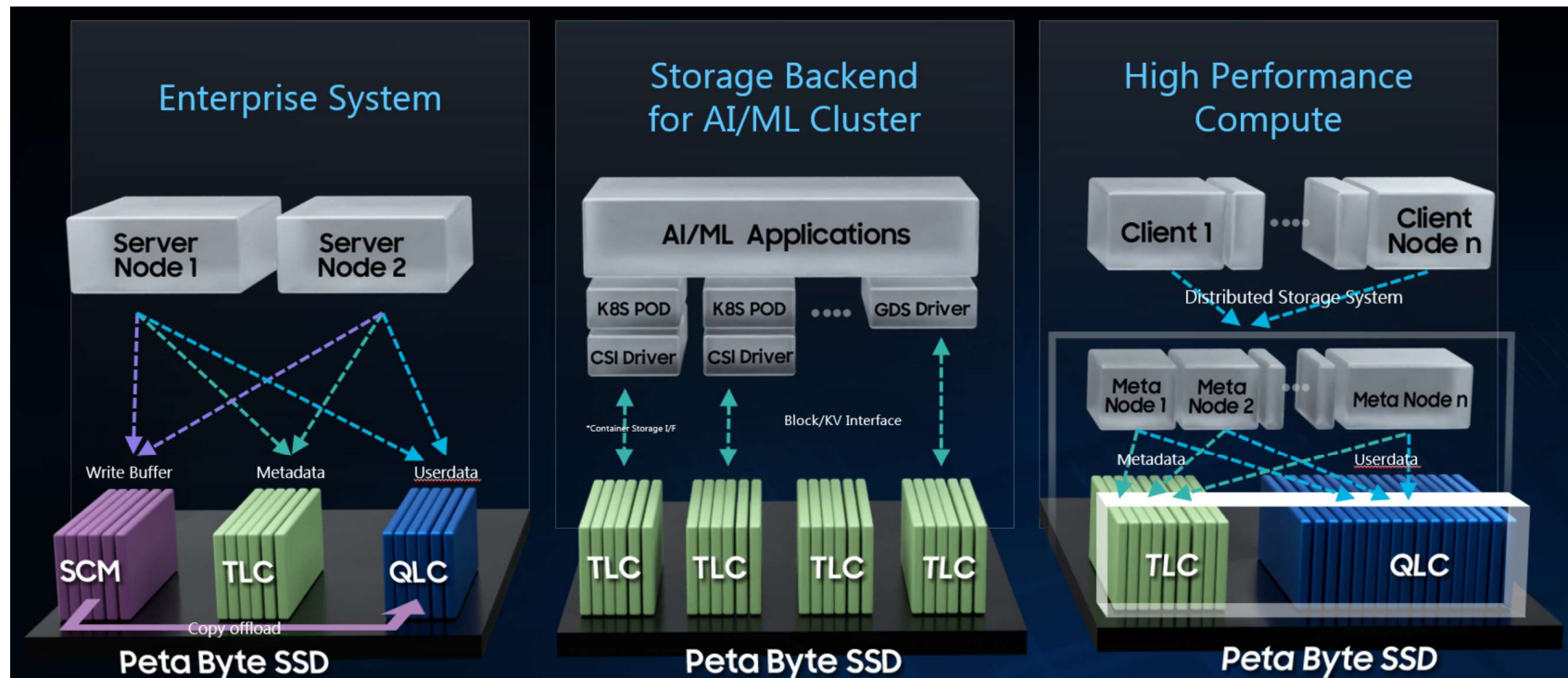
Core Confidential

- Announced at Flash Memory Summit 2022 \* Poseidon V2 is announced at OCP Global Summit 2021



NIC : Dual 200GbE  
TLC : PM1743 PCIe Gen5 V6 TLC 16TB x 8ea (**128TB**)  
QLC : BM1731a V5 QLC 128TB x 16ea (**2PB**)

- Various models can be configured according to the SSD and S/W used.



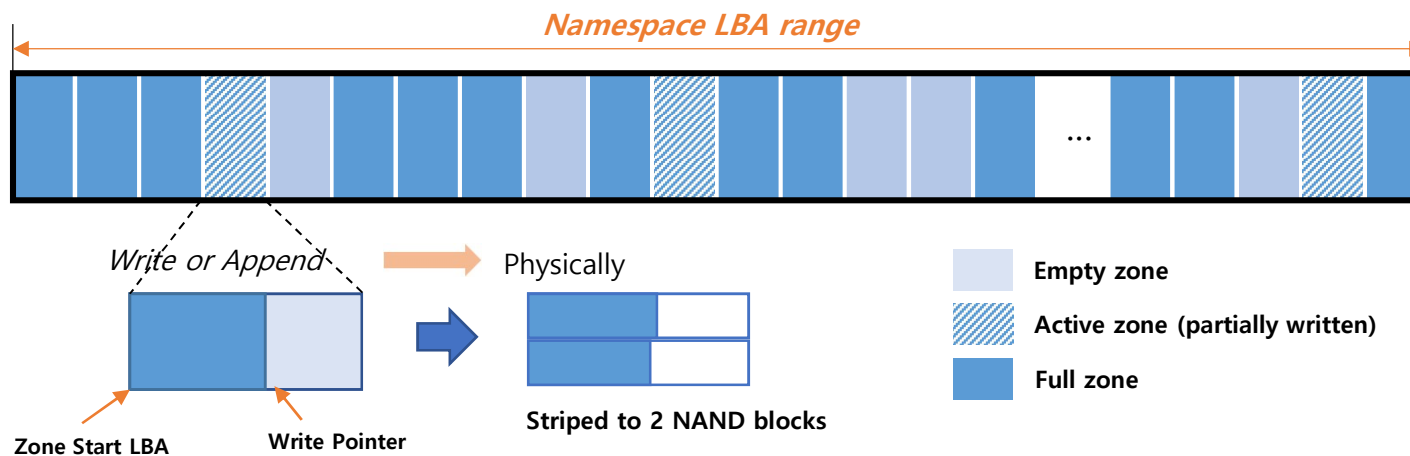


# ZNS Concept (Overview)

Core Confidential

## ■ Zoned Namespace

- ✓ SSD storage space is divided into a fixed size called ZONE ( Zone =? NAND Block Size )
- ✓ Each zone allows only Sequential Write

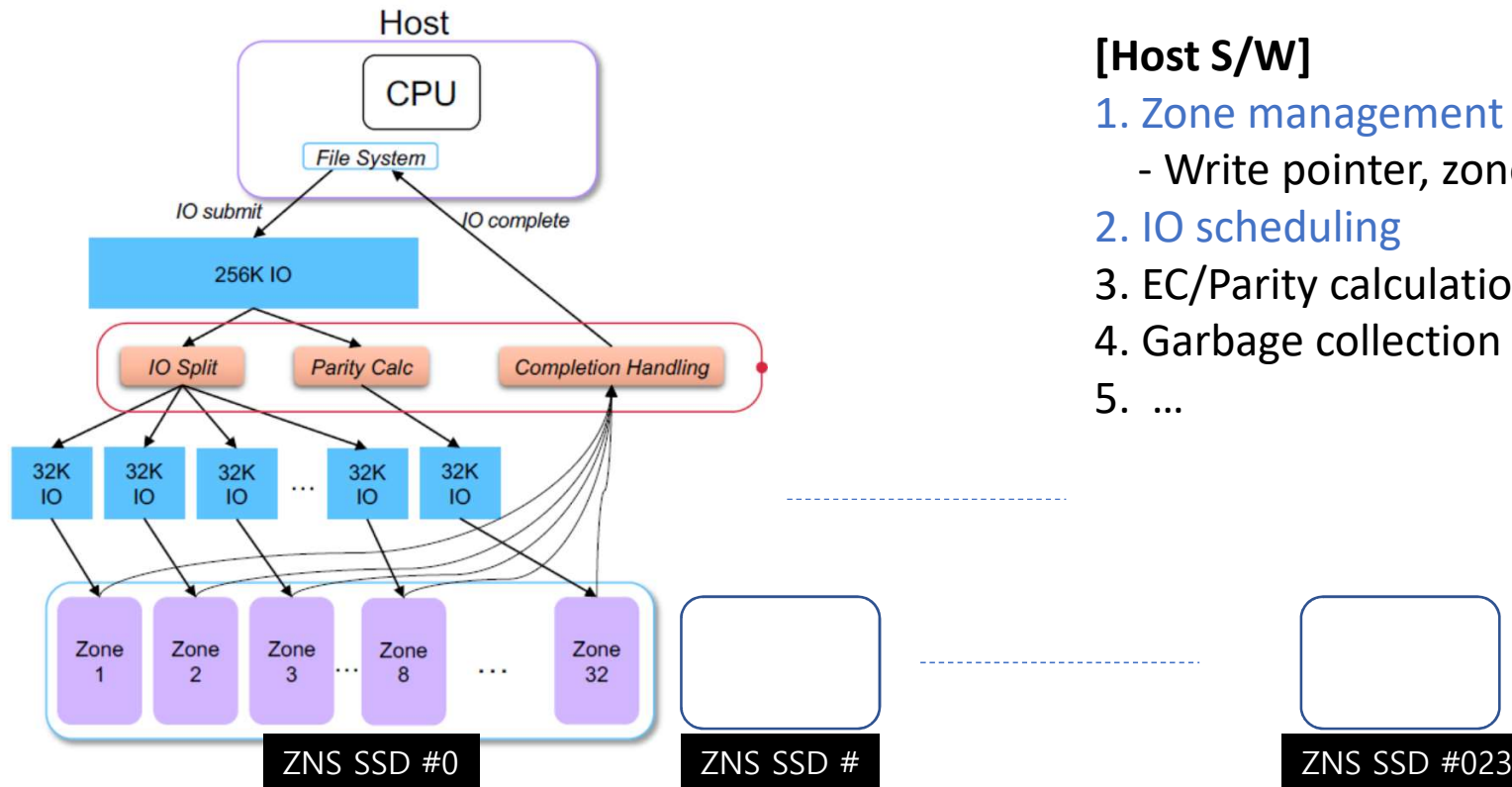


- ✓ Define zone information through zone descriptor and communicate with host

# Challenge 1: Zone & Data Management

Core Confidential

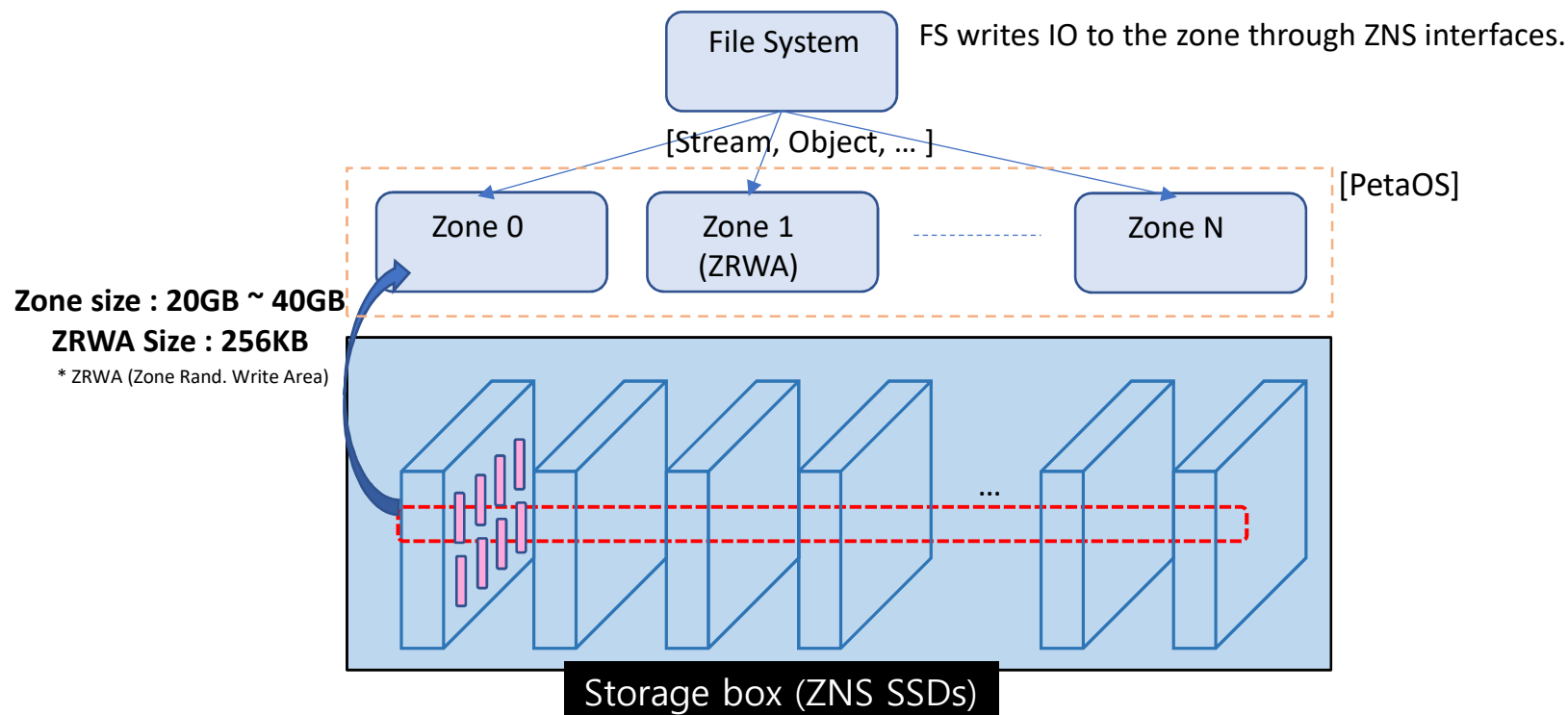
- To use ZNS SSDs, host S/W needs CPU cost to manage zones.



## [Host S/W]

1. Zone management
  - Write pointer, zone activities, ...
2. IO scheduling
3. EC/Parity calculations.
4. Garbage collection
5. ...

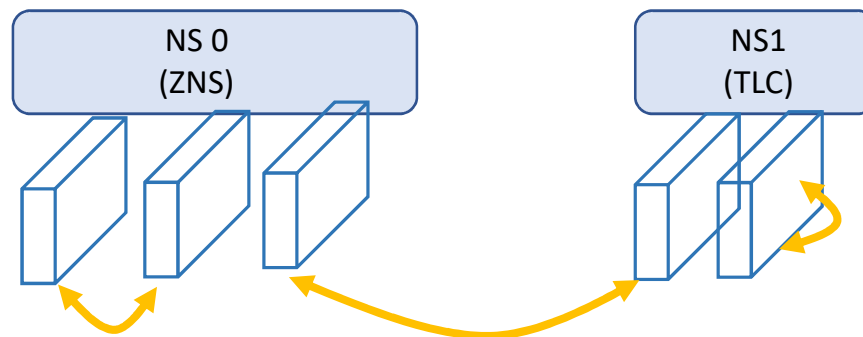
- By offloading device (ZONE) management, it is possible to focus on the original function of host software.



# Challenge 1 : Zone & Data Management

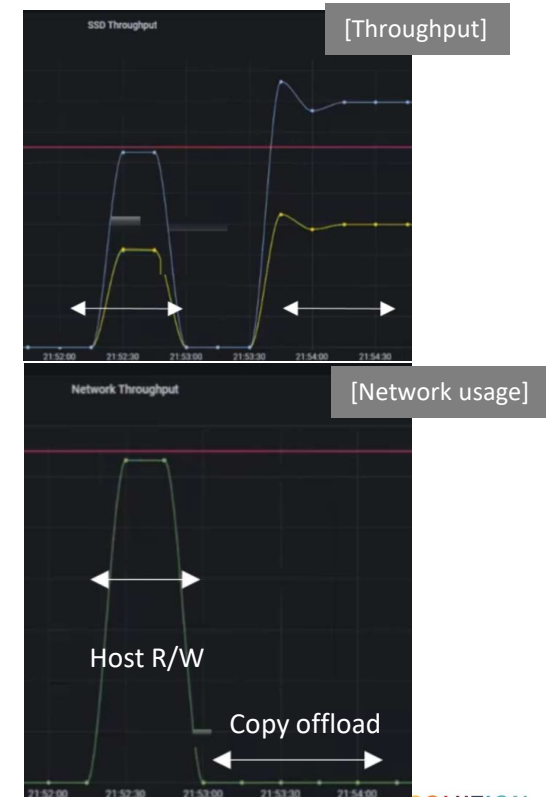
Core Confidential

- PetaOS supports data movement between namespaces
  - Intra/Inter Copy : **Performance can exceed network max**
    - Intra copy : data move inside the NS
    - Inter copy : data move between the NSs



## [Use case]

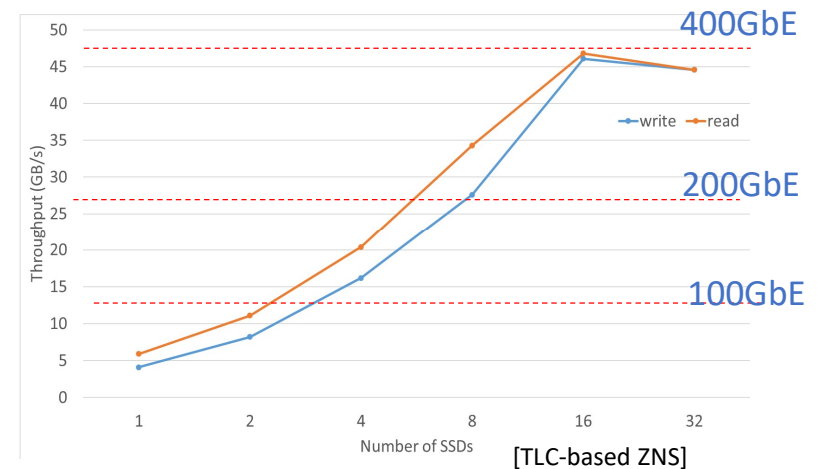
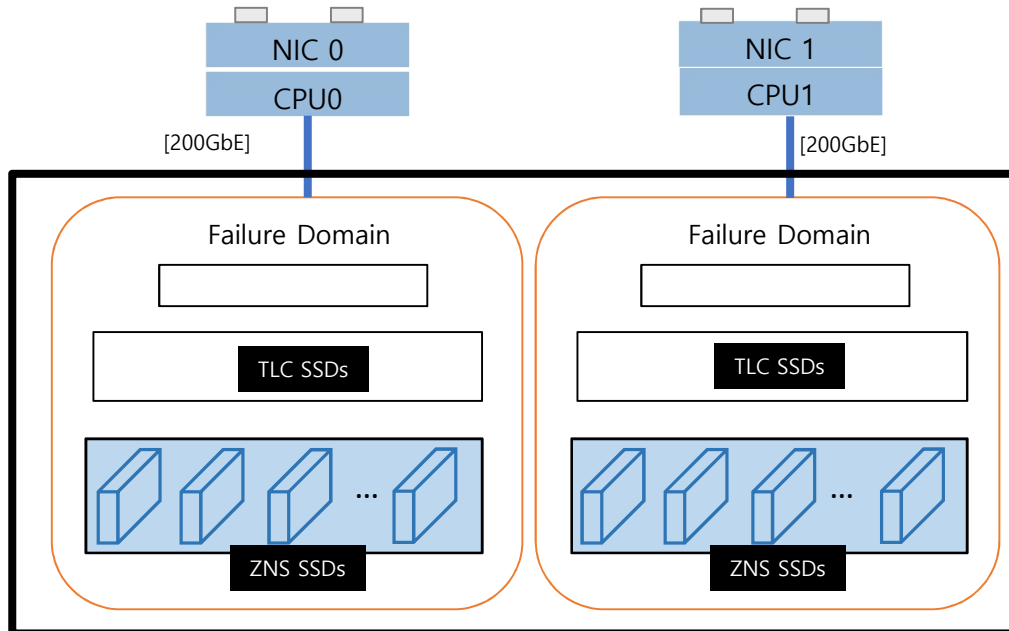
- Control tiering with different type's storages.
- Replication, GC, ...



## Challenge 2 : Blast Radius

Core Confidential

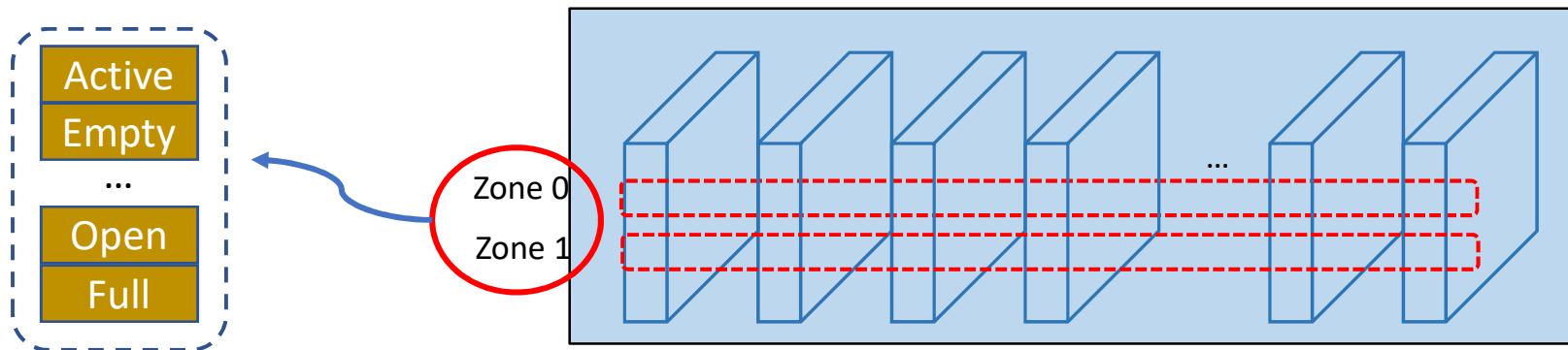
- **Failure Domain (FD) : Isolate to minimize blast radius due to component failure.**
  - How many SSDs are in a FD?



## Challenge 3 : Initialization Time

Core Confidential

- Shortening the init. time of the PB scale system will be competitive.
  - $128\text{TB} * 24\text{ea} = 3\text{PB}$
  - Bring-up time : Server H/W init. + Device map load + PetaOS initial.



**Use zone characteristics as meta info.**

- Save state info. in the CNS => 4sec
- Reconfigure last zone states using limited devices. => 40sec

- **Cheaper Architecture vs. High Performance**

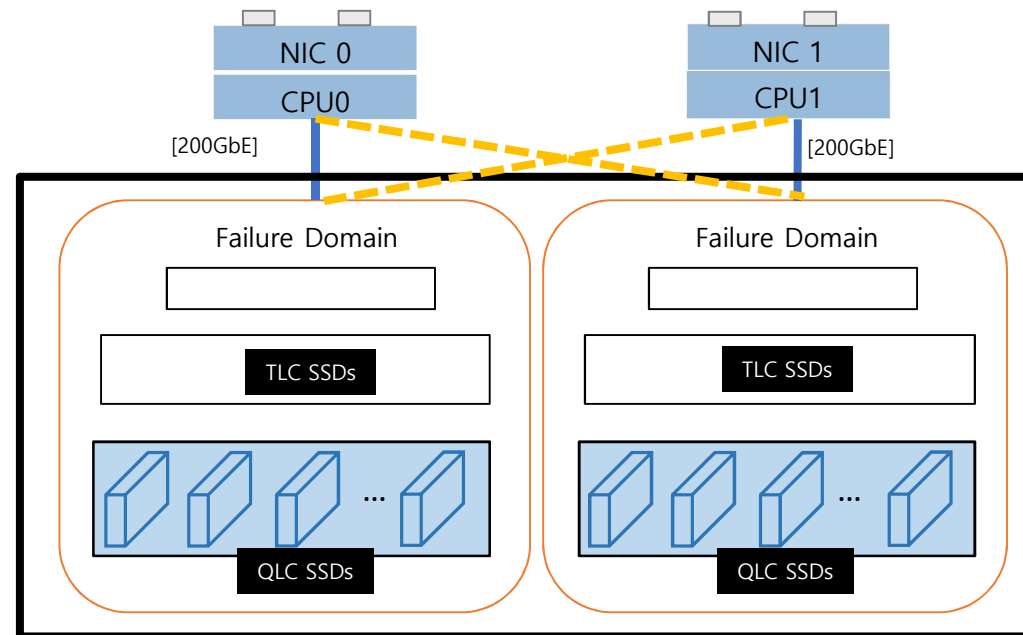
- CPU sockets, DRAM, DPU, ...

- **High-Availability**

- HW & S/W Co-design
- In-box vs. In-cluster (multi-nodes)

- **General vs. Customized System**

- w/ host file system



- **Questions:**
  - How we can reduce the TCO for NAND-based systems?
  - How the host can reduce the management overhead?
- **Our new storage solution**
  - Increasing storage density to Petabyte scale
  - Offload the host management parts into the box.



# SOLUTION

C O R E V A L U E S



**S**peciality  
**O**wnership  
**L**eadership  
**U**pgrowth  
**T**ogether  
**I**ntegrity  
**O**penness  
**N**ow

# Thank You